# Being Actively Ethical: Dynamic UX for AI

Carol J. Smith
Sr. Research Scientist - Human-Machine Interaction, CMU's SEI
Adjunct Instructor, CMU's Human-Computer Interaction Institute

Twitter: @carologic     @sei_etc

Software Engineering Institute
Carnegie Mellon University
Pittsburgh, PA  15213

**Carnegie Mellon University**
Software Engineering Institute

# Copyright Statement

# Responsible, Intentional Design

## Just because you can, doesn't mean you should.

# Early, purposeful work

What is the challenge being face?  Is it AI-friendly?

For whom? What are their needs?

What kind of improvements are expected?

What might a machine do better or faster?

What is not going to be improved (out of scope)?

# AI is a partner - augmenting our abilities

Speed

- Find patterns
- Calculations

Safety (robotics)*

- Dull
- Dirty
- Dangerous
- Dear

*Marr, B. "The 4 Ds Of Robotization: Dull, Dirty, Dangerous And Dear." Forbes. Oct 16, 2017.
https://www.forbes.com/sites/bernardmarr/2017/10/16/the-4-ds-of-robotization-dull-dirty-dangerous-and-dear/#70749ed83e0d

# Diverse teams

Gender, race, culture

Education (school, program, etc.)
Experiences
Thinking process,
Disability status,
and more…

# Not lowering bar ———— extending it

# Diverse, talented and multi-disciplinary

**Includes skill set and problem framing approach**

UX Professionals (big umbrella)

Data Scientists, Machine learning experts

Programmers, System architects

Product managers, etc.

Representatively diverse leadership for retention

Inclusive – Individuals' differences are acknowledged and accepted

# Great Minds Think Different

# High value in diverse teams

Focus more on facts

Process facts more carefully

More innovative

"…become more aware of their own potential [biases](#)"

**Carnegie Mellon University**
Software Engineering Institute

Why Diverse Teams Are Smarter. Harvard Business Review.
https://hbr.org/2016/11/why-diverse-teams-are-smarter

[DISTR]IBUTION STATEMENT A] This material has been approved for public
[releas]e and unlimited distribution.  Please see Copyright notice for non-US
[Govern]ment use and distribution.

10

**Harvard Business Review**

# Ethics for Technology

# AI has great potential, develop with caution

Future AI's *may* be trusted to substitute human cognition and abilities.

Humans must continue to be responsible for situations that involve a person's:

- Life (the use of force)
- Quality of life
- Health
- Reputation

"AI will ensure appropriate human judgement and not replace it" - DIB

# To be biased, is to be human



Bias are shortcuts, to avoid risk and simplify problems.

Not inherently bad,
may be misapplied

Implicit = invisible

Not necessarily in sync
with our conscious beliefs

**Can be managed and changed**

Talk about biases in non-threatening, productive ways

# Biased due to…

Social class

Resource availability

Education

Race, gender, sexuality

Culture, theology, tradition

More…

# All systems have some form of bias

Complete objectivity is misleading.

Bias can have purpose and can be helpful.

The goal is to reduce unintended and/or harmful bias.

# Adopt Technology Ethics

- Harmonize cultural variations
- Balance to pace of change, industry pressure
- Explicit permission to consider and question breadth of implications

# Coalesce on Shared Set of Technology Ethics

Montréal Declaration
Responsible AI_

An initiative of Université de Montréal

1. Well-being
2. Respect for autonomy
3. Protection of privacy and intimacy
4. Solidarity
5. Democratic participation
6. Equity
7. Diversity inclusion
8. Prudence
9. Responsibility
10. Sustainable development

# Diverse, inclusive leaders

# Diverse, Multi-Disciplinary Teams

# Shared Tech Ethics

UX Framework
# Designing Trustworthy AI

**Carnegie Mellon University**
Software Engineering Institute

# Activate curiosity

UX research methods to activate curiosity:

- Abusability Testing
- "Black Mirror" Episodes (inspired by British dystopian sci-fi tv series of same name)
- Flip it to test it
- Implicit Association Test from Harvard

Speculate about system misuse and abuse

- What are potential unintended/unwanted consequences?

More methods to "Outsmart Your Own Biases.": https://hbr.org/2015/05/outsmart-your-own-biases
Implicit Association Test (IAT): https://implicit.harvard.edu/implicit/takeatest.html

# How do we get there?

Montréal Declaration Responsible AI_

An initiative of Université de Montréal

➡ **?** ➡ Trustable, Ethical AI

# Conversations for Understanding

UX Framework guides AI teams

Difficult Topics
- What do we value?
- Who could be hurt?
- What lines won't our AI cross?
- How are we shifting power?*
- How will we track our progress?

*"Don't ask if artificial intelligence is good or fair, ask how it shifts power." Pratyusha Kalluri.
https://www.nature.com/articles/d41586-020-02003-2

**Carnegie Mellon University**
Software Engineering Institute

*"How is this ML model shifting power?" @riakall #NeurIPS2019

# New uncomfortable work

# "*Be uncomfortable*"

## - Laura Kalbag

### Ethical design is not superficial.

# Prompt conversations

**Pair Checklist with Technical Ethics**

- Bridges gap between "do no harm" and reality

**Reduce risk and unwanted bias**

**Support inspection and mitigation planning**

**Carnegie Mellon University**
Software Engineering Institute

Checklist and Agreement - Downloadable PDF:
https://resources.sei.cmu.edu/library/asset-view.cfm?assetid=636620

naterial has been approved for public
se see Copyright notice for non-US

24

# Prompts help reveal hidden tasks

**We work to speculatively identify the full range of risks and benefits:**

☐ Harmful, malicious use and consequences, as well as good, beneficial use and consequences

☐ We will be cognizant and exhaustively research unintended consequences.

**We value honesty and usability:**

☐ Humans can easily discern when they are interacting with the AI system vs. a human.

☐ Humans can easily discern when and why the AI system is taking action and/or making decisions.

☐ Improvements will be made regularly to meet human needs and technical standards.

Checklist and Agreement - Downloadable PDF at SEI:
https://resources.sei.cmu.edu/library/asset-view.cfm?assetid=636620

# UX Framework for Designing Trustworthy AI



1.Accountable to humans

1.Honest and usable

Ethical AI

1.Cognizant of speculative risks and benefits

Respectful and secure

# RightStaff Scenario

AI shift scheduling system

Users: Store managers of fast food restaurants

Goals of RightStaff:

- Faster staffing decisions and scheduling
- Reduced bias of shift selection

# Accountable to Humans

Ensure humans have ultimate control

- Able to monitor and control risk

Human responsibility for final decisions

- Person's life
- Quality of life
- Health
- Reputation

"Ensure humans can unplug the machines"

– Grady Booch

TED Talk, Grady Booch, Scientist, Philosopher, IBM'er

https://www.ted.com/talks/grady_booch_don_t_fear_superintelligence

# Significant decisions

Significant decisions made by the AI system will be

- explained

- able to be overridden

- appealable and reversible

**RightStaff**

- Manager able to reschedule people as needed

# Responsibilities explicitly defined

Between AI system and human(s)

**RightStaff**  (AI System *or Manager?*)

• Picks employees to schedule?

• Defines shifts?

• Method to integrate new information?

  • Sick time

  • Resignations

# Abusability Testing

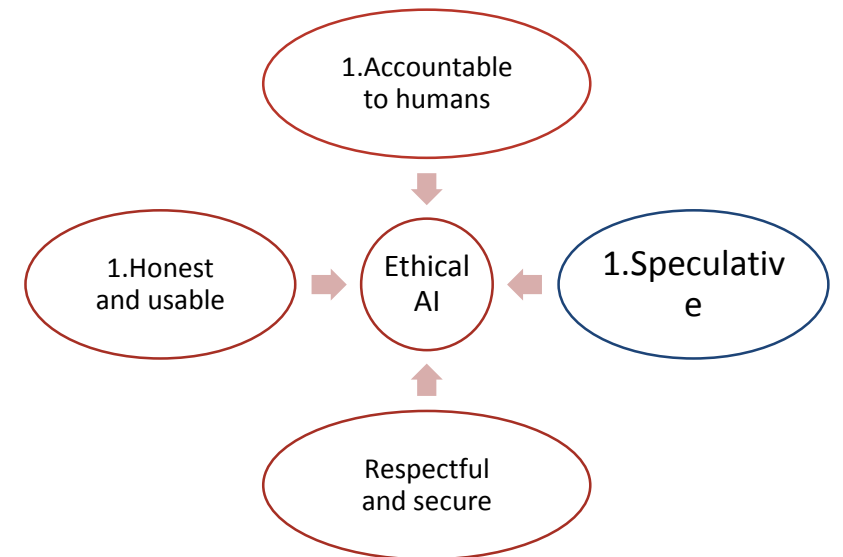Feature added to enable RightStaff to turn off by itself

- What are limits to functionality?

- How could this be abused/misused?

- Implications?

- Risks?

# Cognizant of Speculative Risks and Benefits

Identify full range of

- Harmful, malicious use, as well as good, beneficial use
- Blind spots and unwanted/unintended consequences

# Speculative: Conduct UX research and activate curiosity

Speculate about misuse
and abuse

Potential severe abuse
and consequences

Perspective of people
in frequently marginalized groups

"Black Mirror" episodes

# "Black Mirror" episode

- RightStaff begins prioritizing people with easier schedules

- Managers approve these schedules, reinforcing bias

- People who were previously discriminated against are *still* discriminated against

- What else?

# Speculative: Create communication & mitigation plans

Plan for unwanted consequences

Misuse and abuse of AI system
- Who can report?
- To whom?
- Turn off?
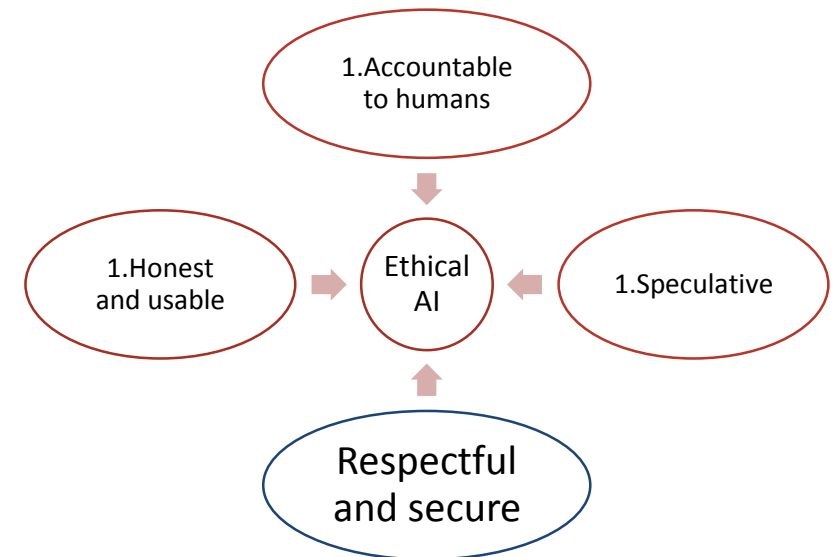- Who notified?
- Consequences?

# Respectful and Secure

Values of humanity, ethics, equity, fairness, accessibility, diversity and inclusion

Respect privacy and data rights

Make system robust, valid and reliable

Provide understandable security

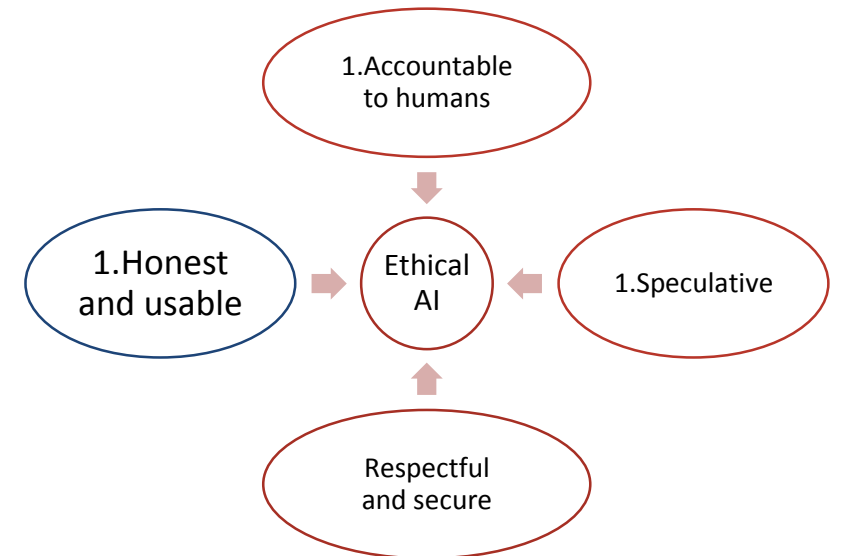# Respectful and Secure

**RightStaff**

- Who has visibility to reasons for changing schedules?

- How is that information used?

- How is PII* of employees protected?

*PII is Personally Identifiable Information (social security number, address, etc.)

# Honest and Usable

Value transparency with the goal of engendering trust

Explicitly state identity as an AI system

# Fair: Remove unwanted bias in data

Show awareness of known and desirable bias

Acknowledge issues

Overcommunicate on issues

**RightStaff**

- System built to reduce the known bias in existing data

- Make it easy to report bias (or prevent it)

# Reward team members for finding ethics bugs

Dr. Ayanna Howard
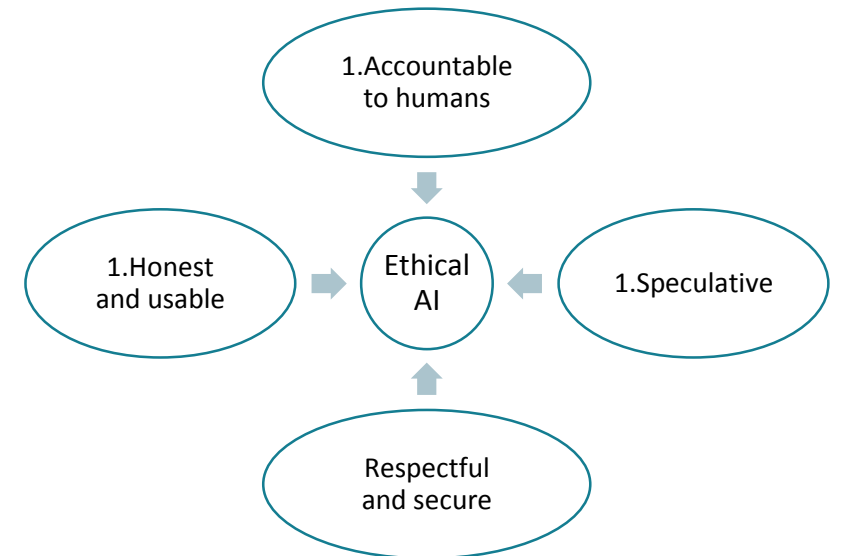- on the Artificial Intelligence Podcast with Lex Fridman

# We aren't perfect, AI won't be perfect

Empower diverse teams, inclusive environments

Adopt technical ethics

Encourage deep conversations (Checklist)

Activate curiosity; be speculative; imaginative

1.Accountable to humans

1.Honest and usable → Ethical AI ← 1.Speculative

Respectful and secure

# Evangelize for human values

# Ethical. Transparent. Fair.

# Carol J. Smith

Twitter: @carologic

LinkedIn: https://www.linkedin.com/in/**caroljsmith**/

CMU's Software Engineering Institute,
Emerging Technology Center

Twitter: @sei_etc